

1/1 WPIL - (C) Thomson Derwent- image  
AN - 1997-303342 [28]  
XP - N1997-250923  
TI - Index generation system for searching document database - associates  
contents of first and second memory using identifier and word  
extracted from partial document as key  
DC - T01  
PA - (XERF ) FUJI XEROX CO LTD  
NP - 1  
NC - 1  
PN - JP09114856 A 19970502 DW1997-28 G06F-017/30 6p \*  
AP: 1995JP-0290408 19951012  
PR - 1995JP-0290408 19951012  
IC - G06F-017/30 G06F-017/21 G06F-017/27  
AB - JP09114856 A

The system has a divider (3) which divides the input document into  
some divisions. An identifier provision part provides the identifier  
providing unit corresponding to each divided document part. A first  
memory (4) matches and stores the identifier and position information  
of each partial document part.

- A word is extracted from each partial document part by an extraction  
part (6). A second memory (7) matches and stores the extracted word  
with the corresponding identifier of each partial document part. An  
index generation unit (8) generates an index used as a keyword, by  
associating the information stored in the first and second memory.
- ADVANTAGE - Enables reduction of index data size by avoiding  
overlapping and describing position information of same word, thus  
reducing memory capacity. Improves document search processing speed.  
(Dwg.1/5)

MC - EPI: T01-J05B3 T01-J11D  
UP - 1997-28

Search statement 32

?ST Y

Session finished: 22 JUL 2002 Time 07:12:44

WPIL - Time in minutes : 6,57  
The cost estimation below is based on Questel's  
standard price list

Estimated cost : ~~17,85 USD~~  
Records displayed and billed : 33  
Estimated cost : ~~49,47 USD~~  
Cost estimated for the last database search : ~~27,02 USD~~  
Estimated total session cost : ~~15,34 USD~~

QUESTEL - Time in minutes : 0,07  
The cost estimation below is based on Questel's  
standard price list

Estimated cost : 0.06 USD

(19)日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11)特許出願公開番号

特開平9-114856

(43)公開日 平成9年(1997)5月2日

(51)Int.Cl. <sup>8</sup>	識別記号	庁内整理番号	F I	技術表示箇所
G 0 6 F	17/30		G 0 6 F 15/413	3 1 0 B
	17/27		15/20	5 5 0 F
	17/21			5 7 0 N
				5 9 0 E
			15/40	3 7 0 A

審査請求 未請求 請求項の数 2 F D (全 6 頁)

(21)出願番号 特願平7-290408

(22)出願日 平成7年(1995)10月12日

(71)出願人 000005496

富士ゼロックス株式会社

東京都港区赤坂二丁目17番22号

(72)発明者 山浦 富久美

神奈川県足柄上郡中井町境430 グリーン

テクなかい 富士ゼロックス株式会社内

(72)発明者 館野 昌一

神奈川県足柄上郡中井町境430 グリーン

テクなかい 富士ゼロックス株式会社内

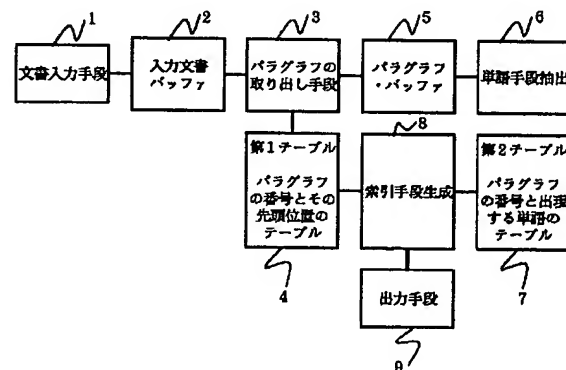
(74)代理人 弁理士 守山 辰雄

(54)【発明の名称】 検索用索引生成装置

(57)【要約】

【課題】 文字数(バイト数)で表現される部分文書の位置情報を同一の語について重複して記述することを回避して、従来に比してデータ量を大幅に減少させた索引を生成する。

【解決手段】 パラグラフ取出し手段3によって入力された文書を複数の部分文書に分割するとともにこれら部分文書に一意の識別子を付与し、第1テーブル4に部分文書の識別子と当該部分文書の文書中における位置情報とを対応付けて格納する。一方、単語抽出手段6によって部分文書から抽出された語を、当該語を抽出した部分文書の識別子と対応付けて第2テーブル7に格納する。そして、索引生成手段8が第1テーブル4と第2テーブル7との格納情報を部分文書の識別子によって関連付けて、語をキーとした索引を生成する。



## 【特許請求の範囲】

【請求項1】 文書中における検索対象の語の位置を、当該文書を構成する部分文書の位置により記述した索引を生成する装置において、  
 文書を複数の部分文書に分割する分割手段と、  
 分割された部分文書に一意の識別子を付与する識別子付与手段と、  
 部分文書の識別子と当該部分文書の文書中における位置情報とを対応付けて格納する第1記憶手段と、  
 部分文書から語を抽出する抽出手段と、  
 語を抽出した部分文書の識別子と抽出された語とを対応付けて格納する第2記憶手段と、  
 第1記憶手段と第2記憶手段との格納情報を部分文書の識別子によって関連付けて語をキーとした索引を生成する索引生成手段と、  
 を備えたことを特徴とする検索用索引生成装置。

【請求項2】 請求項1に記載の検索用索引生成装置において、  
 分割手段は予め設定された基準に従って文書を分割することを特徴とする検索用索引生成装置。

## 【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、電子文書データベースを検索するに際して用いられる索引を生成する装置に関し、特に、全文データベースの索引を生成するに適した装置に関する。

【0002】

【従来の技術】学術文献、書籍、新聞等の電子文書に関するデータベースには、著者や標題等の書誌的事項だけを収録した索引型データベース、書誌的事項に抄録を加えた抄録型データベース、本文自体を収録した全文データベース等、種々の形式があるが、近年では、情報量の豊富な全文データベースが広く利用させている（「情報処理」1992年4月（Vol. 33, No. 4）第413頁～第420頁、「情報処理」1992年10月（Vol. 33, No. 10）第1144頁～第1153頁）。

【0003】データベースに収録されている情報を検索するに際して、当該情報に含まれる語をキーとした索引を利用することにより、検索処理を迅速に行うことができる。特に、全文データベースにおいては、実用的な処理時間で検索を実現するために、索引を用いて検索することが必須とも言える。一般的に、索引には、文書中に出現する語（文字列）と、文書を構成している部分文書の当該語を含むものの文書中における位置が記述されており、この部分文書の位置として、1つの文書を対象とする場合には文書の先頭文字から当該部分文書の先頭文字までの文字数（バイト数）、複数の文書を対象とする場合にはこれら文書識別子と文書の先頭文字から当該部分文書の先頭文字までの文字数（バイト数）との組が用

いられている。なお、部分文書とは、文書を句点、段落等で区切った文書の構成部分をいう。

【0004】図5には、従来の索引の一例を示してある。この索引は、同図（A）に示すテーブルと同図

（B）に示すテーブルとから成っている。（A）に示すテーブルには、或る文書から抽出した語50（出現する語）と、当該語50を（B）のテーブルに関連付けるポインタ51とが含まれており、（B）に示すテーブルには前記語50を含む部分文書の文書中における位置52が含まれている。語の位置52には各語毎の文書中における出現位置をまとめて記述してあり、これら位置52は文書の先頭文字からこれら語の先頭文字までの文字数（バイト数）でそれぞれ表現されている。例えば、“春”という語は、文書中の56バイト目から始まる部分文書、文書中の120バイト目から始まる部分文書、・・・に位置していることが記述されている。

【0005】

【発明が解決しようとする課題】上記したように、従来の索引にあっては、文書中から抽出された多数の語に対してそれぞれ部分文書の位置を記述し、これら部分文書の位置を文書の先頭からの文字数（バイト数）で表現していた。すなわち、同一の部分文書中に含まれる互いに異なる語についてもそれぞれ同一の部分文書の位置情報が格納され、これらの位置情報が文字数（バイト数）で表現されていた。このため、索引全体がかなり大きなデータ量のものとなってしまう、索引を格納するためのメモリの大型化によるコスト増大を招くばかりか、索引に基づく検索処理の遅延化も招くという問題があった。

【0006】本発明は上記従来の事情に鑑みなされたもので、文字数（バイト数）で表現される部分文書の位置情報を同一の語について重複して記述することを回避して、従来に比してデータ量を大幅に減少させた索引を生成する装置を提供することを目的とする。また、これによって索引を小型化し、コスト低減及び検索処理の迅速化を実現する索引生成装置を提供することを目的とする。

【0007】

【課題を解決するための手段】上記目的を達成するため、本発明の検索用索引生成装置では、文字数（バイト数）で表現される位置情報に比べて小さなデータ量となる識別子で部分文書を記述し、文書から抽出した各語にはそれぞれ部分文書の識別子を対応付ける。そして、これら部分文書の識別子には文字数（バイト数）で表現される部分文書の位置情報を対応付け、これによって、部分文書の位置情報を同一の語に対して重複して記述することを回避する。そして、各識別子は各部分文書に一意に対応していることから、各識別子に部分文書の位置情報が重複して対応付けられることもなく、結果として、索引のデータ量が従来に比して大幅に減少する。

【0008】すなわち、分割手段によって文書を複数の

部分文書に分割して、識別子付与手段によって分割された部分文書に一意の識別子を付与し、第1記憶手段に部分文書の識別子と当該部分文書の文書中における位置情報とを対応付けて格納する。一方、抽出手段によって部分文書から抽出された語を、当該語を抽出した部分文書の識別子と対応付けて第2記憶手段に格納する。そして、索引生成手段が第1記憶手段と第2記憶手段との格納情報を部分文書の識別子によって関連付けて、語をキーとした索引を生成する。

【0009】なお、文書を部分文書に分割する位置は、構造化文書におけるタグのように文書中に明示的に記しておいてもよいが、そういった記号がなくても、文書のパラグラフ（節、段落等）を単位とする、文字列の或る長さ（Nバイト）以内のもっとも長い文章の切れ目（句点）を単位とする、文字列の或る長さ（Nバイト）以上の最も短い文章の切れ目を単位とする、といったように種々設定することができる。

【0010】

【発明の実施の形態】本発明に係る検索用索引生成装置を実施する場合の一形態を図面を参照して説明する。図1に示すように、検索用索引生成装置は、索引化する文書データを入力するための文書入力手段1と、入力された文書データを一時記憶する入力文書バッファ2と、一時記憶された文書データをパラグラフ毎の部分文書データに分割するとともに各部分文書データに一意のパラグラフ番号（識別子）を付与するパラグラフ取出し手段3と、付与されたパラグラフ番号と当該部分文書データの位置情報とを対応付けて格納する第1テーブル4と、分割された部分文書データを一時記憶するパラグラフバッファ5と、一時記憶された部分文書データからキーとなる語を抽出する単語抽出手段6と、抽出された語と当該語を抽出した部分文書データのパラグラフ番号とを対応付けて格納する第2テーブル7と、第1テーブル4と第2テーブル7との格納情報から抽出された語をキーとした索引を生成する索引生成手段8と、生成された索引を二次記憶装置等に出力して格納する索引出力手段9と、を備えている。

【0011】文字列入力手段1は全文データベースに収録する文書を検索用索引生成装置に読み込むための手段であり、通常テキストデータとして与えられる文書データを入力文書バッファ2内に格納する。本実施例では部分文書の単位を1パラグラフとしており、パラグラフ取出し手段3は入力文書バッファ2内に記憶された文書データをパラグラフ毎の部分文書データに分割し、各部分文書データを順次パラグラフバッファ5に格納する。

【0012】また、この分割処理に際して、パラグラフ取出し手段3は各部分文書データに一意の識別子を付与するものであり、本実施例では部分文書が文書中の先頭から数えて何パラグラフ目かを示すパラグラフ番号を各部分文書データに付与する。更に、この分割処理に際し

て、パラグラフ取出し手段3は各部分文書の先頭の文字が文書の先頭の文字から数えて何文字目かをカウントし、カウントされた文字数（バイト数）で表現される各部分文書の位置情報を検出する。

【0013】第1テーブル4は読み出し書き込み自在なメモリから構成されており、パラグラフ取出し手段3によって得られた各部分文書毎のパラグラフ番号と位置情報（バイト数）とを対応付けて記憶する。例えば、図3の（A）に示すように、パラグラフ番号“1”の部分文書はバイト数“0”の位置から始まり、パラグラフ番号“2”の部分文書はバイト数“56”の位置から始まり、パラグラフ番号“102”の部分文書はバイト数“86020”の位置から始まるといったように、各部分文書毎のパラグラフ番号と位置情報とを対応付けて記憶する。

【0014】単語抽出手段6は、パラグラフバッファ5に一時記憶された部分文書データを形態素解析してキーとなる語（例えば、自立語）を抽出し、抽出した語と当該部分文書データのパラグラフ番号とを対応付けて第2テーブル7に格納する。なお、部分文書データから語を抽出するためには、形態素解析以外（例えば、DPマッチング法等）の手法を用いることもできる。第2テーブル7は読み出し書き込み自在なメモリから構成されており、抽出された語とパラグラフ番号とを対応付けて記憶する。例えば、図3の（B）に示すように、パラグラフ番号“1”の部分文書から“古来”、“日本人”、“四季”、“変化”、“生活”、“一部”、“ある”という語が抽出された場合には、これらの語をパラグラフ番号“1”で示されるメモリ領域にまとめて格納する。なお、第2テーブル7は第1テーブル4と別途のメモリ装置から構成してもよいが、同一のメモリ装置に領域を分割して第1テーブル4とともに構成してもよい。

【0015】索引生成手段8は、上記のように第1テーブル4と第2テーブル7とに格納されら情報をパラグラフ番号によって互に関連付け、これによって、単語抽出手段6によって抽出された語をキーとした索引を生成する。すなわち、索引生成手段8は、第2テーブル7の格納情報に基づいて、図4の（A）に示すように抽出された語40毎にまとめたテーブル11を作成するとともに、同図の（B）に示すように各語40に対応しているパラグラフ番号42を各語毎にまとめたテーブル12を作成し、テーブル11の各エントリ40（各語）をそれぞれポインタ41でテーブル12の各エントリ42（パラグラフ番号群）に対応付け、更に、テーブル12のパラグラフ番号群42を第1テーブル4の格納情報（同図の（C））に対応付ける。

【0016】上記構成の検索用索引生成装置によると、従来に比してデータ量が大幅に減少した索引が以下のようして生成される。まず、図2に示すテキスト文書が文書入力手段1から入力されると、この文書データがバ

ッファ2に格納される。そして、バッファ2に格納された文書データに対して、パラグラフ取出し手段3が分割処理を行って、部分文書データを順次取り出し、更に、各部分文書データにパラグラフ番号を付与するとともに各部分文書データの位置情報(バイト数)をカウントする。

【0017】これら各部分文書データのパラグラフ番号と位置情報(バイト数)は第1テーブル4に格納される一方、各部分文書データはパラグラフバッファ5に一時記憶されて、当該部分文書データから単語抽出手段6によってキーとなる語が抽出される。そして、これら抽出された語は当該語を抽出した部分文書のパラグラフ番号とともに第2テーブル7に格納され、第1テーブル4と第2テーブル7との格納情報に基づいて索引生成手段8によって図4に示す索引が生成される。

【0018】すなわち、パラグラフ取出し手段3によって、入力文書の1番目のパラグラフ「古来より日本人にとって四季の変化は生活の一部であった。」が取り出されて、バッファ5に格納されるとともに、このパラグラフ番号"1"の部分文書の先頭位置は文書中の0バイト目なので、図3の(A)に示すように、第1テーブル4のパラグラフ番号"1"のエントリには0が書き込まれる。そして、単語抽出手段6によって、バッファ5に格納された部分文書から「古来」、「日本人」、「四季」、「変化」、「生活」、「一部」、「ある」の自立語が抽出され、これらの語が第2テーブル7にパラグラフ番号"1"の部分文書に出現する語として登録される。

【0019】上記と同様にして、入力文書中の2番目のパラグラフ「春はあけぼの。夏は・・・・」が取り出され、パラグラフ番号"2"の部分文書の先頭位置は文書の先頭から56バイト目なので、第1テーブル4の2番目のエントリには56が書き込まれる。このように、第1テーブル4のN番目のエントリの値はN番目のパラグラフの文書の先頭からのバイト数を表している。そして、2番目のパラグラフから「春」、「あけぼの」、「夏」、「夜」・・・・といった語が抽出され、パラグラフ番号"2"の部分文書における出現語として第2テーブル7に登録される。入力文書中の3番目以降の各部分文書についても同様な処理が繰り返され、第1テーブル4及び第2テーブル7に所定の情報が登録される。

【0020】このようにして第1テーブル4と第2テーブル7とが作成された後、索引生成手段8が第2テーブル7の語を重複を排して文字コード順に並べ換え、図4の(A)に示す形式のテーブル11を生成する。また、テーブル11の各語を抽出した部分文書のパラグラフ番号を各語毎にまとめて図4の(B)に示すテーブル12を生成し、テーブル11の各語をポインタ41で関連付けるとともに、テーブル12の各パラグラフ番号を第1

テーブル4のパラグラフ番号に対応付ける。すなわち、索引生成手段8によって図4の(A)、(B)、(C)に示す各テーブル11、12、4がパラグラフ番号で関連付けられ、語をキーとした1つの索引が生成される。

【0021】このように生成された索引は、図5に示した従来の索引に比べて、データ量(バイト数)の多くなる部分文書の位置情報53が重複していない形式となっており、従来に比してデータ量の少ないコンパクトなものとなっている。更に、同一パラグラフ内にある同じ語(例えば、パラグラフ番号"2"の部分文書中にある「春」)は、パラグラフ番号を1つテーブル12に記憶すればよいので、テーブル12のエントリ数を削減することができ、更に索引のデータ量が少なくて済むようになっている。

【0022】上記のようにして生成された索引を用いて検索を行うときには、与えられた検索語によってテーブル11を検索して該当する語を探し、該当する語のエントリのポインタ41が示すテーブル12のエントリを参照して、その検索語が現れる部分文書(パラグラフ)の先頭文字の文書中における位置(なお、次のエントリの値との差から部分文書の長さも)を得ることができる。そして、検索結果として、該当する部分文書の先頭部分等をディスプレイ装置等に表示する。

【0023】なお、検索語はの部分文書中における出現位置は部分文書の位置として得られ、検索語の詳しい位置がわからないので、語と語の位置関係(例えば、「特許」と「出願」が隣り合って現れる、「特許」の後20文字以内に「出願」が現れる等)といった検索は索引を用いた処理だけでは行うことができない。しかしながら、索引から検索語の有無及びその検索語が含まれる部分文書が特定できるので、その部分文書のテキストそのものを参照して検索語と検索語との位置関係を確認すれば、所期の目的を容易に達成することができる。また、部分文書を1パラグラフや1ページといったような小さな単位に設定すれば、そのテキストを参照して検索語同士の位置関係を確認する作業は、あまり時間をとらずに容易に行うことができる。

【0024】なお、部分文書の単位の設定に、パラグラフ、セクション、節、章といった階層をもたせ、生成された索引の利用範囲を拡大してもよく、このように部分文書単位の階層を表現する場合にあっても、索引全体から見れば容量の小さなテーブル12を拡張するだけであるので、従来に比して索引をコンパクトなものに維持することができる。

【0025】

【実施例】約320Kbyteの入力文書データについて、上記した本発明に係る装置で索引を生成したところ、パラグラフ数は約1000、文書中に出現する異なる語の数は約6700、図4の(B)に示すテーブル12のエントリ数は約27000、1エントリは2byte

eであり、同図の(C)に示すテーブル4のエントリ数は約1000、1エントリは3byteであった。すなわち、これらテーブル11、4の容量はそれぞれ、54Kbyte、3Kbyte程度であった。

【0026】一方、図5に示した従来技術による索引では、上記のテーブル11及び4に該当するテーブル(同図の右側部分)のエントリ数は約59000、1エントリは3byteとなり、当該テーブルの容量は約177Kであった。したがって、従来と同一なテーブル11

(図4の(A))を除いて比べると、従来では177Kbyte必要であったものに対し、本発明によれば約1/3の57Kbyteで済み、索引全体としても大幅にデータ量が削減されたことが確認された。

【0027】

【発明の効果】以上説明したように、本発明に係る検索用索引生成装置によると、文書に含まれる語の位置を特定するための部分文書を一意に付与した識別子によって識別して、バイト数表現される部分文書の位置情報を同一の語に対して重複して記録することを回避するようにしたため、従来に比してデータ量を大幅に減少させて索引を生成することができ、索引の小型化によって必要\*

\*とするメモリ容量を減少させてコスト低減を達成することができ、更には、索引を用いた検索処理の迅速化を達成することができる。

【図面の簡単な説明】

【図1】 本発明の一例に係る検索用索引生成装置を示す構成図である。

【図2】 索引を生成するための入力文書の一例を示す図である。

【図3】 第1テーブル及び第2テーブルに格納される情報を示す概念図である。

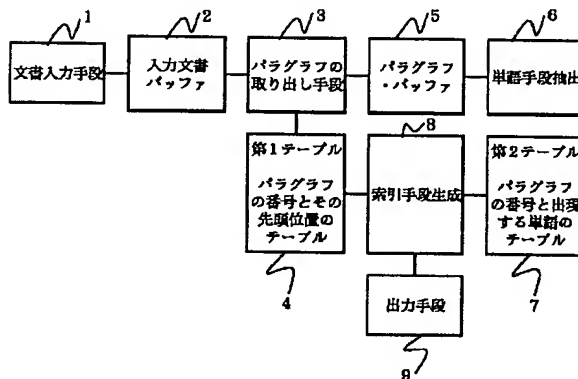
【図4】 本発明の検索用索引生成装置で生成される索引の一例を示す概念図である。

【図5】 従来の検索用索引の一例を示す概念図である。

【符号の説明】

1・・・文書入力手段、 2・・・入力文書バッファ、  
3・・・パラグラフ取出し手段、 4・・・第1テーブル、  
5・・・パラグラフバッファ、 6・・・単語抽出手段、  
7・・・第2テーブル、 8・・・索引生成手段、  
9・・・出力手段、

【図1】



【図2】

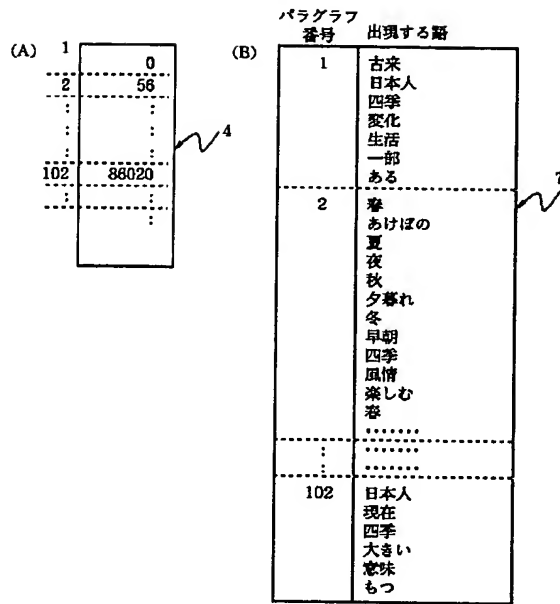
入力文章：

古来より日本人にとって 四季の変化は生活の一部であった。

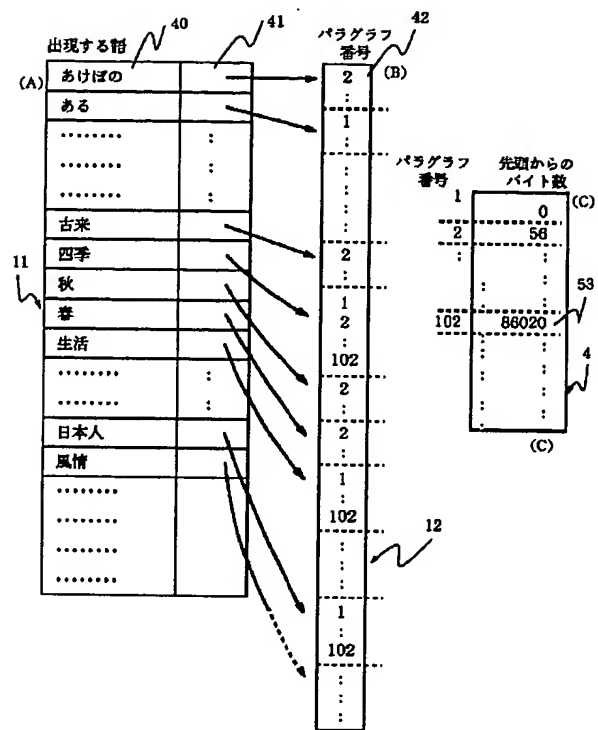
春はあけぼの。夏は夜。秋は夕暮れ。冬は早朝。四季の風情を楽しむ。春の .....

日本人にとって、現在も四季は大きな意味をもっている。 .....

【図3】



【図4】



【図5】

